

Form Filling based on Constraint Solving

Ben Spencer, Michael Benedikt, Pierre Senellart

► To cite this version:

Ben Spencer, Michael Benedikt, Pierre Senellart. Form Filling based on Constraint Solving. 18th International Conference on Web Engineering (ICWE 2018), Jun 2018, Cáceres, Spain. 10.1007/978-3-319-91662-0_20 . hal-01851555

HAL Id: hal-01851555

<https://hal.inria.fr/hal-01851555>

Submitted on 30 Jul 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

Form Filling based on Constraint Solving

Ben Spencer¹, Michael Benedikt¹, and Pierre Senellart²

¹ Department of Computer Science, University of Oxford

² DI ENS, ENS, CNRS, PSL University & Inria & LTCI, Télécom ParisTech

Abstract. We describe a system for analyzing form-based websites to discover sequences of actions and values that result in a valid form submission. Rather than looking at the text or DOM structure of the form, our method is driven by solving constraints involving the underlying client-side JavaScript code. In order to deal with the complexity of client-side code, we adapt a method from program analysis and testing, concolic testing, which mixes concrete code execution, symbolic code tracing, and *constraint solving* to find values that lead to new code paths. While concolic testing is commonly used for detecting bugs in stand-alone code with developer support, we show how it can be applied to the very different problem of filling Web forms. We evaluate our system on a benchmark of both real and synthetic Web forms.

1 Introduction

Finding data on the Web is useful for search, information extraction, and aggregation. The massive scale of the Web and the data itself means these tasks must necessarily be automated. Interesting data is often hidden behind user interfaces, and in particular *Web forms*. For example, airline or real-estate websites provide free access to their data via search forms, but there is no public API or standard way to access this information. This data is “hidden” from Web search engines and other automated tools, and makes up the *deep* or *hidden Web* [8, 21].

An automated tool accessing Web forms must find some sequence of user actions (for example clicking buttons or filling input fields) which leads to a successful form submission: one that leads to the target data. These actions might involve complex interactive user-interface elements such as drop-down lists, date pickers, and tabs. In addition, they must satisfy certain restrictions on the actions and inputs given such as mandatory and optional fields, and validation rules for input values. This combination of actions and input values creates a huge search space for form filling tools to consider. The restrictions on input actions and values are normally enforced in the browser directly with client-side JavaScript (for usability) and again on the server when the query is received (for security). Human-focused interfaces and input validation rules make it difficult for automated tools to correctly fill and submit the forms.

Much research has been devoted to effective, automatic form-filling [5, 25, 26, 30, 32, 45], or even further, to generating *wrappers* which use the form to look up and extract data [18, 40, 43]. A common target is search fields, which generally have no validation constraints [5, 26, 39]. Other work assumes domain knowledge, encoded in heuristics or rules [18]. For complex fields in the absence of domain knowledge, approaches to finding field values include sampling text on the page [28], and the application of machine

Listing 1 Validation code for the example airline form.

```
function validate_to() {  
    var from = document.getElementById("from").value;  
    var to = document.getElementById("to").value;  
    return validate_aux(from, to);  
}  
function validate_aux(from, to) {  
    if (from.length == 0) { alert("Error: Departure airport must be set"); return false; }  
    if (from == to) { alert("Error: Departure must differ from Destination"); return false; }  
    return true;  
}  
function validate_date() {  
    var to = document.getElementById("from").value;  
    var date = document.getElementById("date").value;  
    if (to.length == 0) { alert("Error: Destination Airport must be set"); return false; }  
    var day = parseInt(date.substr(0, 2), 10), month = parseInt(date.substr(3, 5), 10);  
    var now = new Date();  
    var valid = (month >= now.getMonth()+1 && (month != now.getMonth()+1 || day >= now.getDate()));  
    if (!valid) { alert("Error: date cannot be before today"); return false; }  
    return true;  
}
```

learning techniques [25, 34]. Without both domain knowledge and a corpus of examples, these approaches cannot infer constraints enforced by client side code, and thus cannot find satisfying values. We aim to supplement these approaches with *constraint-driven form filling*, which analyzes client-side code to determine constraints being enforced on form actions and values, and solves these constraints to yield successful submissions. Constraint-driven form filling allows form exploration tools to exploit the rapid advances in constraint-solving technology [6, 13, 27].

Example 1. Consider an airport form, which includes fields **From** and **To** (with identifiers “from” and “to”, respectively) with values populated from a drop-down list, along with field **Date** (with identifier “date”) populated by a date picker. A snippet of the form’s validation code is shown in Listing 1, where `validate_to` is attached to the **To** field, and `validate_date` is attached to the **Date** field. For **To**, the code checks that the **From** field has already been filled and that it is not equal to **To**. For **Date** it checks that **To** is already filled and that the date entered is later than today’s date (for simplicity, all dates are assumed to be in the current year).

Note that the validation code involves restrictions on both the values of the fields, and on the order in which they are filled. In addition to the constraints explicitly enforced by the event handlers, there are a number of implicit constraints on the field: for example, **From** and **To** are implemented by drop-down lists, so the values can only be chosen from a certain pre-determined set. A form-filling or wrapper-generation tool must find values that satisfy all of these constraints.

Static analysis of the underlying code can in principle determine the set of restrictions enforced by the code (and thus values leading to form submission). Unfortunately, analyzing JavaScript code used on the Web is notoriously difficult [7]. Thus a popular intermediate position in testing JavaScript code is *dynamic analysis*, with a well-known approach being *concolic testing* [20, 42]. Concolic testing combines concrete and symbolic execution, using the symbolic analysis and constraint solving to generate new concrete input values to test. The goal is to generate test inputs which cover each distinct code path in the program or function being tested.

We introduce FormSolve, which applies the idea of mixing execution, symbolic tracing, and constraint solving to the exploration of Web forms. FormSolve will generate input actions and values for Web forms, symbolically trace the code executed with these inputs, and solve constraints that will lead to new inputs which reach new code paths.

We propose a refinement of concolic testing tuned for Web forms. This involves many challenges that are not encountered in concolic testing of standalone code. Instead of using a constraint solver to find input values for a single function, our concolic testing algorithm explores sequences of user input events, along with their corresponding values, that may trigger a set of event-handling functions. Our browser infrastructure allows the symbolic tracing of this event-handling code, which can be written either directly in pure JavaScript or using popular libraries such as jQuery. It also provides fine control over the browser, which is necessary for faithfully—and deterministically—emulating user form-filling actions, as well as for controlling and monitoring the browser’s behavior. In order to get interesting output, our goal is not to generate arbitrary input values (as in exhaustive testing), or interesting corner cases (as in functional testing), but rather values that a typical user might provide via the interface. Our constraint generation allows us to focus on *user-realizable values*, thus accelerating the discovery of interesting form submissions. We evaluate our technique on both synthetic and real-world forms, and compare it to alternative approaches to Web form filling.

Related work. Indexing and extraction from the deep Web involves a number of tasks, including entry point finding [4, 33, 34], form label identification [17, 37, 46], form-filling [5, 25, 26, 30, 32, 45], and result page analysis [11, 47]. Work on form-filling is focused on finding values that extract a good set of results. For example, some work attempts to choose keywords for a text field which return relevant results [5]. Tools such as Crawljax [36] and AJAX Crawl [14] take into account the state of the user interface while crawling. The emphasis there is in identifying changes in the DOM (*Document Object Model*, the tree structure representing the content and structure of a Web page) caused by user actions and events, not on getting through forms. Our work is on form-filling, but the focus is only on getting results that satisfy validation rules. Additional desired properties can be overlaid on top of this. One of our challenges is choosing an interesting ordering for the form filling events. This challenge also arises in the analysis of other event-driven systems, such as Android apps [1, 22, 24].

Concolic testing is a well-established topic, applied to many programming languages, including JavaScript. Yet there are few attempts to apply it to the Web. SymJS [29] uses concolic testing, but is based on the Rhino JavaScript engine, which can parse only a small fraction of real-world Web JavaScript. Jalangi [41] is a framework allowing instrumentation and runtime monitoring of JavaScript. Neither SymJS nor Jalangi support Web forms. The only systems we know of for exploring forms via JavaScript analysis are the demonstration systems ProFoUnd [7] and ArtForm [44] (precursor to this work).

2 Concolic testing for Web form exploration

Concolic testing, or directed automated random testing [20, 42], is a testing technique which uses concrete executions of a program to drive a symbolic analysis. The symbolic analysis guides the automated tester and suggests specific inputs which reach new parts of the program which would be difficult for traditional test-generation approaches to

Algorithm 1 High-level algorithm for concolic testing.

```
1: procedure CONCOLIC-ANALYSIS(program)
2:   values  $\leftarrow$  CHOOSE-INITIAL-VALUES()
3:   trace  $\leftarrow$  EXECUTE-AND-RECORD(program, values)
4:   path-tree  $\leftarrow$  INIT-TREE(trace)
5:   while path-tree is not fully explored do
6:     target-path  $\leftarrow$  SEARCH(path-tree)
7:     values  $\leftarrow$  SOLVE-PATH-CONSTRAINT(target-path)
8:     if target-path was successfully solved then
9:       trace  $\leftarrow$  EXECUTE-AND-RECORD(program, values)
10:      EXTEND-TREE(path-tree, trace)
11:     else
12:       MARK-UNREACHABLE(path-tree, target-path)
13:     end if
14:   end while
15: end procedure
```

discover. Conversely, the concrete executions allow exploration through parts of the program which are not fully understood by the symbolic analysis alone.

A generic concolic testing algorithm is given in Algorithm 1. Concolic testing begins by choosing some default starting values for the variables. The function is then executed concretely (that is, using a real interpreter) and symbolic information is recorded about how the inputs are modified and when they occur in branch conditions. Thus each trace is associated with a *path condition*: the set of individual branch conditions (expressed in terms of the input values) which must be satisfied or unsatisfied in order for the program execution to follow that particular path. Each path condition is a logical formula describing an equivalence class of input values, with equivalent values resulting in the same execution path in the program. Thus the state of the exploration can be characterized by the set of path conditions of explored traces, which form a tree. The testing algorithm proceeds in an execute-and-analyze loop. After execution of a trace, the new information gathered during that trace is added to the tree. In order to find set of input values that reaches a new path, the algorithm chooses a sequence of branch conditions that has not been explored, generates the corresponding path condition, and sends it to a *constraint solver*. If the solver can satisfy the constraint, a solution is chosen as the next set of input values. If the solver cannot, the path is marked as unreachable. This process is repeated to cover progressively more code on each new iteration.

Example 2. Consider the function `validate_aux(from,to)` from the running example. Concolic testing would first execute the program on random values for the arguments `from` and `to`, say empty strings for both. This would bring the program to a trace that terminates after the first alert. Tracing this path symbolically, the algorithm identifies the constraint `from.length = 0`, which is the first path added to the tree (line 4). The search command on line 6 will isolate the path consisting of the single constraint $\neg(\text{from.length} = 0)$ as an unexplored path in the code, and this will be sent to a constraint solver (line 7), which will return values for `from` and `to` that satisfy it. For example, the solver might return the values `from = 'a'` and `to = 'a'`. These values are used to execute the function in the next iteration, which drives the code to the second alert. This second execution is symbolically traced, with the path $\neg(\text{from.length} = 0) \wedge \text{from} = \text{to}$. The search procedure

on line 6 now selects $\neg(\text{from.length} = 0) \wedge \neg(\text{from} = \text{to})$ as an unexplored path to target. Solving this constraint will give values that drive the execution to avoid each alert.

Note that concolic testing normally includes *classification* of a trace—in the context of testing, this would be determining whether an error occurs in it. In Algorithm 1, we assume this is done within EXECUTE-AND-RECORD. For our form crawling application, the classifier determines whether the result of a trace is a successful submission (leading to a new page) or not. We can detect this by checking for page loads, **alert** messages, or certain modifications to the DOM.

Adaptation to form exploration. In adapting concolic testing to the exploration of Web forms, a tempting analogy is that the form represents a “virtual function”, whose arguments are the form fields. However, a function takes its arguments all at once, whereas a Web form takes inputs one at a time interactively. The code which responds to a form filling can be triggered by a number of events, such as filling fields or hovering over them. These event handlers may interact, and the order in which they are fired can affect the output. Thus a more realistic model of a form is as an association of each field f with one or more actions Act_f , where an action writes the input value v_f into the field and executes some program code. In using concolic testing to explore a form, we must discover both the values to fill into form fields and *also* the ordering of these actions.

Our algorithm is shown in Algorithm 2. We explain the details, and in particular the constraint solving call on line 12, below. It is similar in structure to a classic concolic algorithm, with the high-level distinction that we explore a tree of paths for each action. At each iteration we choose an action with unexplored code, and generate both a value for each form field and an action ordering.

We now formalize the problem of simultaneous detection of input values and the ordering of actions using an extremely idealized model. Although the model’s assumptions do not hold on real Web forms, the approach is still effective.

In the model we have a set of *form actions*, each with a corresponding program which manipulates a set of variables \vec{v} , including one distinguished variable v , the *input variable* of the form action. Informally this is the value supplied by the user when invoking this action. Each variable v has a *default value* D (provided by the Web page’s HTML) which it holds until the action is run. In our basic model, we assume program code in a simple procedural language built up from the following grammar:

$$x := \tau(\vec{y}) \quad \text{if } \varphi(\vec{y}) \text{ then } E_1 \text{ else } E_2 \quad E_1; E_2 \quad \text{Abort} \quad \text{do}_C$$

Above τ ranges over terms built up from set of atomic functions (e.g., $+$, \times) from variables and constants, while φ ranges over some set of conditions (e.g., Boolean combinations of atomic conditions $\tau_1 \{ \leq, \neq, = \} \tau_2$, where τ_i are terms). C in do_C ranges over some set of commands that do not impact submission or control flow. A condition φ is *ground* if it contains no free variables.

A variable is *assigned* in expression E if it occurs on the left side of some assignment statement and otherwise is *free*. The semantics of the language are standard. Given an expression E and a binding σ for the free variables of E , the semantic function returns the sequence of ground conditions and atomic actions **Abort** and do_C that are generated during an execution. We will be particularly interested in the truth values of conditions and whether or not **Abort** is encountered. Given a condition φ and a binding σ for the

Algorithm 2 The high-level algorithm for form exploration.

```

1: procedure FORM-CONCOLIC-TESTING(form-actions)
2:   initial-order, initial-values  $\leftarrow$  CHOOSE-INITIAL-ORDER(form-actions)
3:   trace  $\leftarrow$  EXECUTE-AND-RECORD(initial-order, initial-values)
4:   path-tree  $\leftarrow$  INIT-TREE(trace)
5:   if trace is terminating then
6:     Mark all local symbolic paths within trace as Known-Extendible
7:   else
8:     Mark aborting local symbolic paths as Known-Unextendible
9:   end if
10:  while  $\exists \text{Act} \cdot \text{Unresolved}_{\text{Act}}(\text{path-tree}) \neq \emptyset$  do
11:    target  $\leftarrow$  CHOOSE(sympth, Act) with sympth  $\in$   $\text{Unresolved}_{\text{Act}}(\text{path-tree})$ 
12:    order, values  $\leftarrow$  SOLVE-PATH-CONSTRAINT( $\text{OverApprox}_{\text{Act}}(\text{target})$ )
13:    if target was successfully solved with order, values then
14:      trace  $\leftarrow$  EXECUTE-AND-RECORD(order, values)
15:      EXTEND-TREE(path-tree, trace)
16:      if trace is terminating then
17:        Mark all local symbolic paths within trace as Known-Extendible
18:      else
19:        Mark aborting local symbolic paths as Known-Unextendible
20:      end if
21:    else
22:      Mark target as Known-Unextendible
23:    end if
24:  end while
25: end procedure

```

variables in φ , the *value* of φ , denoted $\text{Val}(\varphi, \sigma)$, is true if φ holds when the variables appearing in φ are replaced with their valuation by σ . The *trace* of E on an assignment σ to the free variables of E is the sequence of conditions encountered and their values. If E executes **Abort** at any point when running on σ it is said to *abort on* σ , while otherwise we say it *terminates on* σ .

We name our actions by numbers, with Act_i denoting the program associated with action i , while v_i and D_i denote the distinguished input variable and default value of Act_i , respectively. A set of such indexed actions $\text{Act}_1 \dots \text{Act}_n$ has *restricted global state* if for every Act_i , each free variable v occurring in its program expression is one of the input variables v_j and no input variable is ever assigned. That is, the actions have no shared global state except the input variables, which are set as each action is executed.

The trace of an action on a binding σ is simply the trace of its program code expression. The behavior of a single expression is well-defined given a binding for all variables. We now need to explain the outcome of a sequence of form actions, which involves binding each free variable v encountered in conditions to either the user-specified input value of v or to its default value.

A *bound form action* is a pairing of a form action Act_i with a value c_i for its input variable v_i , while a *form input* is a sequence of bound form actions. We define the *unfolding* of a form input $(\text{Act}_1, c_1) \dots (\text{Act}_n, c_n)$ as the sequence $(E_1, \sigma^1), \dots, (E_n, \sigma^n)$, where σ^i is the *order-modified assignment* mapping v_j to c_j if $j \leq i$ and to D_j otherwise. We can extend our semantics to form inputs via unfoldings. The *trace* of a form input

is the concatenation of the traces in its unfolding. A form input is said to abort if some E_i aborts on σ_i ; otherwise we say it terminates.

We now extend our discussion from traces to symbolic descriptions of these traces, starting with the “local behavior” within a form action, given a binding for all variables. For binding σ to the free variables of expression E , the *local symbolic path* of $E(\sigma)$ is a formula describing the values of conditions in the trace of $E(\sigma)$. That is, the conjunction $\bigwedge_{\text{Val}(\varphi_i, \sigma) = \top} \varphi_i \wedge \bigwedge_{\text{Val}(\varphi_i, \sigma) = \perp} \neg \varphi_i$ where $\varphi_1 \dots \varphi_k$ is the sequence of conditions encountered in executing E on σ .

We have just discussed the symbolic paths that emerge from concrete traces. We now describe formulas that represent possible concrete traces for future exploration. Given a trace $t = (\varphi_1, \text{TVal}_1) \dots (\varphi_k, \text{TVal}_k)$ for an action, consisting of conditions φ_i and their truth values TVal_i , a *symbolic modification* is a local symbolic path of the form $(\varphi_1, \text{TVal}_1) \dots (\varphi_p, \text{TVal}_p), (\varphi_{p+1}, \neg \text{TVal}_{p+1})$ where $p < k$. That is, we take a prefix of t , and negate its last element. Given a set of traces T , a symbolic modification is *unexplored* if no trace in T satisfies this condition. Given set of traces T and action A , $\text{Unexplored}_A(T)$ is the set of symbolic modifications of traces in A that are unexplored.

We can lift the classification of traces as aborting or terminating to the symbolic level. Given a set of traces T and a form action A , let $T(A)$ be the restriction of the trace to A . Let $\text{Abort}_A(T)$ be the traces in $T(A)$ that are aborting, and $\text{Terminate}_A(T)$ the set of traces that are terminating. By $\text{Symbolic}(\text{Abort}_A(T))$ we denote the set of local symbolic paths of traces in T that are aborting in A , and similarly for $\text{Symbolic}(\text{Terminate}_A(T))$.

Above we have a symbolic version of the behavior of a single form action A , describing it via conditions. We now need to lift this to a sequence of form actions. In doing this we have to take into account the role of the ordering in determining whether we use the default value or the user-supplied value. Further, since we are interested in whether a path is explored in a terminating trace, we also need to track symbolically whether *other* actions abort.

Given a local symbolic path $\psi(\vec{v})$ of a trace for Act_i , the *ordered version* of ψ , $\text{Ord}_i(\psi)$, is the formula ψ with each variable v_j replaced by v'_j and conjoined with the constraint $(j \preceq i \rightarrow v'_j = v_j) \wedge (j \succ i \rightarrow v'_j = D_j)$ using an additional relation \preceq . That is, $\text{Ord}_i(\psi)$ is a symbolic representation of an ordering and a binding σ such that ψ holds on the order-modified assignment σ^i defined above.

Given a path sympth for action Act_i , and a set of traces T , let $\text{OverApprox}^i(\text{sympth})$ to be $\text{Ord}_i(\text{sympth}) \wedge \text{OrderAx} \wedge \bigwedge_{j \neq i} \text{MayTerminate}_j$, where MayTerminate_j is defined as $\bigwedge_{\text{sympth}'_j \in \text{Symbolic}(\text{Abort}_{\text{Act}_j}(T))} \neg \text{Ord}_j(\text{sympth}'_j)$ and OrderAx is a formula asserting that \preceq is a linear order on the indices of actions $1 \dots n$. That is, the formula $\text{OverApprox}^i(\text{sympth})$ symbolically represents the orderings and values that will achieve the behavior sympth in action Act_i and will not drive any other action to a known-aborting trace of T . Informally, it describes form inputs that may explore sympth without aborting in any action, based on the current knowledge of aborts in T .

Similarly, $\text{UnderApprox}^i(\text{sympth})$ is $\text{Ord}_i(\text{sympth}) \wedge \text{OrderAx} \wedge \bigwedge_{j \neq i} \text{MustTerminate}_j$, where MustTerminate_j is $\bigvee_{\text{sympth}'_j \in \text{Symbolic}(\text{Terminate}_{\text{Act}_j}(T))} \text{Ord}_j(\text{sympth}'_j)$. This formula represents orderings and values which achieve the behavior sympth in action Act_i and drive every other action to known-terminating trace of T . Informally, this describes inputs and orderings that we are sure will explore sympth based on what we know about aborts in T . The following two lemmas give the critical properties of these symbolic descriptions.

Lemma 1. *If $\text{UnderApprox}^i(\text{sympth})$ is satisfiable by input values $c_1 \dots c_n$ and ordering $j_1 \dots j_n$, then the form input $(\text{Act}_{j_1}, c_{j_1}) \dots (\text{Act}_{j_n}, c_{j_n})$ gives a trace that does not abort outside of action Act_i and extends sympth .*

Lemma 2. *If sympth is a symbolic modification for action Act_i , and the form input $(\text{Act}_{j_1}, c_{j_1}) \dots (\text{Act}_{j_n}, c_{j_n})$ generates a trace extending sympth that does not abort outside of Act_i , then $c_1 \dots c_n$ and $j_1 \dots j_n$ satisfy $\text{OverApprox}^i(\text{sympth})$.*

Lemmas 1 and 2 form the basis of Algorithm 2. We maintain a set of traces and corresponding symbolic paths which, along with their symbolic modifications, form a tree. We also classify the symbolic paths and their modifications. We distinguish the *known-extendible* paths, those which are known to have an extension that is terminating, and the *known-unextendible* ones, where it is known that there is no such extension. The paths which are neither known-extendible or known-unextendable are said to be *unresolved*. The set $\text{Unresolved}_A(T)$ contains the local symbolic paths from T which are neither known-extendible or known-unextendible, as well as the symbolic modifications of local symbolic paths in T which have not yet been considered by the algorithm.

At any step of the algorithm we choose an unresolved path sympth and check whether $\text{OverApprox}^i(\text{sympth})$ is satisfiable. If the formula is not satisfiable, we mark sympth as known-unextendible (that is, no terminating trace extends sympth , which we know from Lemma 2). Otherwise we take a satisfying assignment consisting of $c_1 \dots c_n$ and ordering $j_1 \dots j_n$, and use it in a new execution, giving trace t . We add the new trace t to our set of traces and iterate.

If t terminates, then it acts as a witness that each restriction of t to action Act_i can be extended by a terminating trace. Thus, we mark all the local symbolic paths in t as known-extendible. If t aborts, then by Lemma 1 we know that either Act_i aborted, or $c_1 \dots c_n$ and $j_1 \dots j_n$ did not satisfy $\text{UnderApprox}^i(\text{sympth})$. In either case, there must be some action Act_j for which t does not extend an explored branch, and where t aborts. That is, in Act_j , t follows a previously unexplored path and discovers an abort. Thus one local symbolic modification is resolved in Act_j (even if it had not been a known modification until now).

It is possible for a newly recorded trace to give rise to *new* symbolic modifications, so the total number of unresolved traces does not necessarily decrease at each iteration. However, each action's program code has a finite number of symbolic branches, so the number of symbolic paths available to explore is limited. Because at least one local symbolic path is resolved in each iteration, $\bigcup_{i \in \{1..n\}} \text{Unresolved}_{\text{Act}_i}(T)$ must eventually become empty, and this guarantees termination.

Proposition 1. *Algorithm 2 is complete. Assuming completeness of the solver, on termination the path tree will have the property that for every local symbolic path sympth which has any extension which terminates overall, then at least one such extension is explored by the algorithm. In other words, every local symbolic path which is reachable on a terminating trace is explored.*

Example 3. Let us illustrate Algorithm 2 on Example 1. Assume that we have identified the form actions as Act_{From} , Act_{To} , and Act_{Date} for entering values into the departure airport, arrival airport, and departure date fields. Algorithm 2 will choose an arbitrary default initial order and values for these fields: for example empty strings for the airports

and “01/01” for the date. The code is executed on these values, and reaches the first alert message in `validate_aux`, and the first alert in `validate_date`, both classified as aborts. The corresponding local symbolic path for Act_{To} is the single constraint $\text{from.length} = 0$, and for Act_{Date} it is $\text{to.length} = 0$. These constraints are added to the corresponding path trees for each action. The command on line 11 will then choose an unexplored path to target, suppose it is $\neg(\text{to.length} = 0)$ in Act_{Date} . Assuming the default value D_{To} is the empty string, the corresponding ordered constraint Ord_{Date} simplifies to: $(\text{Act}_{\text{To}} \prec \text{Act}_{\text{Date}} \wedge \neg(\text{to.length} = 0))$. To form the full over-approximation constraint, we also include the linear order axioms, and `MayTerminate`, which in this case (when simplified) is simply the negation of the single abort trace in Act_{To} : $(\text{Act}_{\text{From}} \prec \text{Act}_{\text{To}} \wedge \neg(\text{from.length} = 0))$.

The call on line 12 will solve this combined constraint for both an order and values. The only valid order is $\text{From} \prec \text{To} \prec \text{Date}$. Suppose the returned values are $\text{To}_1 = \text{“A”}$, $\text{From}_1 = \text{“A”}$, and $\text{Date}_1 = \text{“01/01”}$. The code is re-tested with these values, and symbolically traced, leading to the second alerts in both `validate_aux` and `validate_date`. This second trace is associated with local symbolic paths for each action, for example $\neg(\text{from.length} = 0) \wedge \text{from} = \text{to}$ for Act_{To} , which are added to the corresponding trees on line 15. In the second iteration of the **while** loop we would choose another unexplored path (line 11), and this would return, for example, the path $\neg(\text{from.length} = 0) \wedge \neg(\text{from} = \text{to})$ from Act_{To} . The ordered version of this path simplifies to: $(\text{Act}_{\text{From}} \prec \text{Act}_{\text{To}} \wedge \neg(\text{from.length} = 0) \wedge \neg(\text{from} = \text{to}))$. This time, the `MayTerminate` constraint is required to avoid both aborts in Act_{Date} . Thus, the full over-approximation constraint, omitting the linear order constraints, simplifies to:

$$\begin{aligned} & [\text{Act}_{\text{From}} \prec \text{Act}_{\text{To}} \wedge \neg(\text{from.length} = 0) \wedge \neg(\text{from} = \text{to})] \\ & \wedge [\text{Act}_{\text{To}} \prec \text{Act}_{\text{Date}} \wedge \neg(\text{to.length} = 0)] \\ & \wedge \left[\text{int}(\text{substr}(\text{date}, 3, 5)) \geq m_1 \wedge \left(\bigvee \begin{array}{l} \text{int}(\text{substr}(\text{date}, 3, 5)) \neq m_1 \\ \text{int}(\text{substr}(\text{date}, 0, 2)) \geq d_1 \end{array} \right) \right] \end{aligned}$$

where m_1 and d_1 are the concrete date and month observed during the execution.

A second call to the solver returns the ordering $\text{From} \prec \text{To} \prec \text{Date}$ and values From_2 , To_2 , Date_2 , where From_2 and To_2 are distinct and non-empty, and Date_2 is later than the current date. These new values are tested, this time producing a trace which terminates in every action and successfully submits the form.

The completeness result is based on strong assumptions. It requires completeness of the solver, and requires the analysis to track all conditions symbolically. The code must also conform to the simple structure where form actions set their input variables but otherwise do not communicate with other actions. Real-world code does not obey these assumptions, but we can still apply Algorithm 2 to arbitrary code, tracking only the input fields symbolically across actions, while dropping the completeness guarantees.

3 Implementing concolic testing for form exploration

Implementing concolic testing in the Web context is challenging. First, one needs to control the browser, both to simulate user actions faithfully and ensure that the browser behaves deterministically from one iteration to the next. Frameworks such as Selenium WebDriver support this, but give limited control over certain low-level events such as

timers and AJAX events. Second, we need to get information from the browser, and in particular need to analyze the executed code to record the symbolic paths taken. There are frameworks which can be used for symbolic tracing of stand-alone JavaScript, with the goal to help testing and debugging [41]. However, applying these to third-party JavaScript on the Web can be problematic, since they require instrumentation of the JavaScript source, which can be problematic when crawling third-party sites.

FormSolve’s approach works directly with a browser engine. We build on top of WebKit, an open-source production Web browser engine, used by Apple’s Safari browser. The browser engine includes page fetching, HTML and CSS rendering, and a JavaScript interpreter, but excludes the browser’s user-interface. We instrument the source code to add hooks to control the browser, and we instrument the interpreter to perform symbolic-tracing, extending the WebKit interpreter to track a symbolic value with every concrete value. WebKit’s JavaScript interpreter uses an internal JavaScript bytecode language, so our symbolic tracing works at the *byte-code level*, rather than the JavaScript source level. This simplifies the generated constraints.

The symbolic interpreter runs alongside WebKit’s existing JavaScript interpreter and computes the symbolic values used in the analysis. When values are read from form fields, they are tagged as symbolic, and these symbolic values are propagated by the symbolic interpreter as the values are processed. A branch instruction is called symbolic if its branch condition uses any symbolic value. When a symbolic branch is encountered, its condition is recorded in the trace.

Symbolic values are created when certain properties of DOM objects are accessed by the site’s JavaScript code. For example the `value` property of text fields, the `checked` property of checkboxes and radio buttons, or the `value` or `selectedIndex` properties of drop-down boxes. To make these DOM properties symbolic, WebKit’s internal implementation of the DOM API is instrumented.

As well as the bytecode interpreter, JavaScript’s built-in methods must also be instrumented. These are implemented by C++ methods in WebKit, which are modified in FormSolve to add an appropriate symbolic tag to the return value. We have only instrumented the functions which were most commonly used on the sites we were analysing. This reduces the implementation complexity, but also simplifies the generated constraints. Branch conditions which cannot easily be encoded as SMT constraints are dropped to make the analysis more concrete, thus allowing the analysis to continue past certain difficult functions or patterns.

Note that JavaScript library functions (such as from jQuery) are not built-ins and are implemented in JavaScript, so they are executed along with any other JavaScript code. They do not require special handling by the interpreter.

Example 4. We explain our symbolic tracing on this snippet from Example 1:

```
var from = document.getElementById("From").value;
if (from.length == 0) { alert("Error: Departure Airport must be set"); }
```

Table 1 shows the corresponding bytecode, the resulting register-level state change, and the changes in the symbolic state generated by the symbolic interpreter. For example, the first line of the table shows that the initial JavaScript command generates two bytecode instructions: **op_call** `r1 getElementById "from"`, which calls `getElementById` to look up the DOM node with identifier “from” (the `From` field) and store it in the register `r1`, followed by **op_get_by_id** `r2 r1 "value"`, which fetches the “value” property of that DOM

Table 1. JavaScript level, bytecode-level, and symbolic execution.

JavaScript code	WebKit bytecode	Concrete state	Symbolic state
<code>var from = document. getElementById("from").value;</code>	<code>op_call r1 getElementById "from"</code> <code>op_get_by_id r2 r1 "value"</code>	$r1 := \text{DOM node From}$ $r2 := \text{"" [value of From]}$	(none) $r2 := \text{SymStr("from")}$
<code>if (from.length == 0)</code>	<code>op_get_by_id r3 r2 "length"</code> <code>op_eq r4 r3 0</code> <code>op_jfalse r4 else_label</code>	$r3 := 0$ [length of ""] $r4 := \text{true}$ [length is 0] [$r4$ is true ; so no jump]	$r3 := \text{StrLen}(r2)$ $r4 := \text{IntOp}(r3, \text{LEQ}, \text{ConstInt}(0))$ $b := \text{BoolOp}(r4, \text{B_EQ}, \text{ConstBool}(\text{false}))$

object (the empty string) and stores it in register $r2$. The symbolic instrumentation sets $r2$ to the symbolic value of the property lookup, which will be **SymStr**("from"), representing a symbolic input originating from a field with identifier "from".

Finding user-realizable values. Line 12 of Algorithm 2 uses the constraint solver to generate new input values. In the context of form filling, we want to avoid inputs that a user could not perform at the interface, since these are unlikely to produce a useful output. We thus add additional constraints on the value space that enforce *user realizability*.

HTML permits various types of input fields, each with their own restrictions on which values can be input. For example a *select* element produces a drop-down list with fixed options for the user to select from. To produce user-realizable input values, the analysis models these input fields to only generate values a user would be able to provide using a normal Web browser. This is done by encoding DOM facts as extra constraints which are included with each path condition.

Example 5. In the running example, the `To` field is implemented by a drop-down to choose between a set of airport codes: JFK, ORD, etc. In any constraint involving the corresponding variable `to`, we add a constraint saying $\text{to} = \text{JFK} \vee \text{to} = \text{ORD} \vee \dots$, where the list of codes is populated from the DOM. The client-side code may also check the index of the selected item, rather than its value. In this case, we add a variable `to_index` to represent this index, which must correspond to the main `to` variable. In our example, the constraint would become $(\text{to} = \text{JFK} \wedge \text{to_index} = 1) \vee (\text{to} = \text{ORD} \wedge \text{to_index} = 2) \vee \dots$, forcing the solver to choose a matching index and value together.

Constraint solving. A critical component in the architecture is the constraint solver. In general, concolic testing requires a high-performance solver which supports standard program variable types: integers, Booleans, reals, bit-vectors, and arrays. Examples of solvers used for concolic testing or symbolic execution include Boolector [38], STP [19], Yices [15], and Z3 [12]. JavaScript requires very strong support for reasoning about strings. In particular, transformations between strings and other data types are common, so it is very useful if the solver supports these datatype coercions. We make use of *CVC4*, which supports a wide variety of theories compared with other solvers [9], strong string support [31], as well as support for coercions between different theories.

4 Experimental Evaluation

Standalone JavaScript code. The first goal of our evaluation is to compare our bytecode-based code exploration with alternatives based on JavaScript source. We compared against Jalangi [23, 41], a framework for JavaScript testing that includes concolic testing

Table 2. Comparing FormSolve with Jalangi on 1000 synthetic programs.

	FormSolve Jalangi	
Avg. distinct paths	5.75	2.54
Avg. line coverage	91.49%	87.69%
No. with full coverage	46.44%	37.05%

Table 3. Results of FormSolve, Crawljax and Artemis on 1000 synthetic forms.

	FS	CJ	Ar
Analysis time (s)	969 272 153	2 206	
Avg. iterations	1.86	40.95	43.29
No. with full cov.	10.4%	—	0.7%
Forms submitted	31.4%	19.8%	17.7%

Table 4. Results of FormSolve, Crawljax and Artemis on 18 JSFiddle examples.

	FS	CJ	Ar
Analysis time (s)	2 245	6 509	319
Avg. iterations	26.17	33.67	45.94
Forms submitted	61.1%	33.3%	11.1%

Table 5. Comparing Alg. 2 with single-order testing on 1000 synthetic forms.

	Single-order Alg. 2	
Analysis time (s)	585	969
Avg. iterations	3.04	1.86
No. with full coverage	4.7%	10.4%
Forms submitted	24.1%	31.4%
Multiple traces explored	69.8%	46.8%
↳ average line coverage	90.6%	95.5%
↳ no. with full coverage	6.7%	22.2%
↳ forms submitted	34.0%	66.2%

for standalone (that is, not Web-based) JavaScript. We compared the tools on 1000 synthetic, randomly-generated standalone JavaScript programs. To generate the examples, we began with a context-free grammar for the program expressions, and then used a standard CFG-sampling algorithm [35] to generate random expressions from it. We then choose a random implementation for each operation in the generated program skeleton, choosing from a fixed set of simple JavaScript implementations. For example, for a term $a \wedge b$, we choose between a single-expression implementation combining the two child expressions, one which uses intermediate variables, and one using conditionals.

The aggregate results are shown in Table 2. We see that FormSolve explores many more paths than Jalangi. FormSolve benefits from symbolic instrumentation of a number of built-ins, native string operations, and datatype coercions which are not supported (and may be difficult to support) at the source level in Jalangi. There is also some benefit from FormSolve using CVC4, over the earlier version of CVC used by Jalangi.

Synthetic Web forms. The second goal was to test FormSolve’s ability to deal with complex Web forms. We ran FormSolve on a suite of 1000 randomly generated test forms. The generator first produces a random form, and then supplements it with validation code using the JavaScript generator described above. Each field’s validation code uses the value of its own field, and optionally those of other fields. The generated forms use an average of 8.57 form fields and 47.16 lines of JavaScript validation code. We compared FormSolve with two alternative approaches. Crawljax [10, 36] is an automatic crawler for dynamic websites which is based on a dynamic analysis of the Web application. In particular, Crawljax tracks changes to the state of the DOM, and explores user events until no more states can be discovered. In the absence of any domain knowledge, Crawljax can be seen as a state-of-the art approach to crawling complex websites, including Web forms. Artemis [2, 3] is a tool that does feedback-directing testing of websites. In each iteration, Artemis generates DOM and JavaScript events as well as values to enter in

form fields. It uses metrics like line coverage to rank actions that are most promising for use in the next iteration. Each tool (FormSolve, Crawljax, Artemis) is run for a maximum of 50 *iterations*: for FormSolve, an iteration is defined as in Alg. 2; for Crawljax, it is a full run of the tool, until no new state is discovered—different iterations use different random seeds; for Artemis, an iteration is a predefined sequence of actions that is run in the browser. Iterations in Crawljax take a large amount of time and involve many browser interactions, while they are comparatively short in FormSolve and Artemis.

The results are shown in Table 3. The table includes line coverage of the event-handling code, measured at the JavaScript source level, for all tools except Crawljax, which does not record this. The number of iterations for each example is taken to be that of the first submission, or the total if no submission was found. While not truly representative of the kinds of constraints found on real websites, these results do indicate that FormSolve can deal with complex constraints that do *not* obey the restricted global state assumptions our completeness result relies on, while also showing that prior methods do not suffice for complex constraints.

Real-world form validation. A more interesting benchmark is given by a set of simple real-world examples of Web forms. We used JSFiddle, a code-sharing website for Web-based code. Since form-related examples are not easily isolated in JSFiddle, we performed a Google search for “JSFiddle form validation”. We extracted 18 examples, after removing those which did not perform JavaScript form validation. The examples were modified where necessary to allow them to be run in the tools with common submission and error behaviors. These test cases represent real-world forms; they use common JavaScript libraries and use complex constraints, but do not include the full complexity of real sites. There is no other content to contend with than the forms themselves.

We again compared FormSolve with Crawljax and Artemis, with results shown in Table 4. Some examples required only non-empty inputs; Artemis still has difficulty handling some of these, since it does not faithfully simulate the form-filling events triggered by a real user. FormSolve has a modest gain over Crawljax, handling several of the examples that require string comparisons.

The running time for FormSolve is dominated by a single example which uses string functions to parse an ID number, including a checksum; which took 30 minutes 45 seconds out of the total analysis time for all 18 examples of 37 minutes 25 seconds. This example (named *VSKNx* in our test data [16]) requires solving many difficult constraints over both strings and integers. At one point characters from the ID are converted to integers, multiplied by a constant, converted back to strings to extract the digits, which are then converted to integers and summed; and the analysis must solve a final integer constraint on the cumulative total of these sums. As such, FormSolve does not find a submission within 50 iterations. If left to run with no iteration limit, it is eventually able to find a valid ID number and successfully submit the form. In total, there were three examples which ran to the iteration limit, all of which can be solved with enough time. Including those examples would bring FormSolve’s successful submissions to 77.8%.

Impact of dynamic re-ordering. We compared Algorithm 2 to a variant of the algorithm where a fixed field order (taken from the field ordering in the DOM) is used. The experiments were conducted on the same set of 1000 synthetic forms described above. The JSFiddle examples are not tested, because they do not include interesting interactions

between different fields’ event handlers. For each configuration, we measured the line coverage as before, and also tracked whether we were able to generate a successful form submission. Note that as the examples are randomly generated, some of them have no valid inputs, and submission really is impossible. The results are shown in Table 5.

We observe that the big advantage of Algorithm 2 over a static ordering is pruning out explorations in one event handler which are known to abort in a later handler. The algorithm is able to prove that certain unexplored paths must later lead to abort actions, and therefore avoid wasted iterations. Trace execution is very expensive when testing real sites, so this is a significant advantage.

An example of this can be seen in the second part of the table. Algorithm 2 only attempted to explore 46.8% of the examples, compared to 69.8% for the static-ordering variant. In the cases explored by the static ordering algorithm but rejected by Algorithm 2 (23.0% of the total), the static ordering algorithm spent time exploring branches in certain event handlers which would all eventually lead to an abort action in a later handler. In contrast, Algorithm 2’s calls to the constraint solver showed these branches to be futile, allowing the algorithm to ignore them. This explains the dramatic difference in the number of successful form submissions found in examples where multiple traces were explored (34.0% for the static order and 66.2% for Algorithm 2).

Real websites. Finally we experimented with 4 real airline websites that included complex validation rules. We report the numbers for the single-order mode discussed above. In that mode, FormSolve was able to find correct submissions while Crawljax could not (despite taking more than 16 hours in total) and Artemis could find only 2. An illuminating example is the archived website of AirTran Airways. The validation rules include that origin and destination must be entered, that there must be valid combinations of the number of children and adults (e.g., not more infants than adults), and that the departure and return dates must be ordered correctly and be no earlier than today’s date. FormSolve finds 10 successes while exploring 207 code paths. Handling of the date and passenger constraints exercises the support for integer arithmetic in the solver. The Australian airline Rex has similar passenger restrictions, but in addition each departure airport allows a different set of arrival airports. FormSolve can find 458 distinct successful paths to submission, of a total of 1539 paths explored.

The real-world set is only anecdotal, due to our prototype’s basis in the development version of WebKit, which limits its ability to run on recently-updated websites. Further, the dynamic reordering-mode has not yet been made sufficiently robust—for example, it has weaker support for restricting to user-realizable values. This causes it to miss half of the submissions found by the single-order mode on our small set of real-world sites. More thorough experiments with a more modern browser codebase will be needed to draw firm quantitative conclusions about the general applicability of our approach.

5 Conclusion

In this paper we make the first step at applying constraint-solving technology to finding valid submissions for Web forms. The approach has the advantage that it can be used not just to find single submissions but a representation of all valid submissions, in terms of a constraint; this is particularly relevant for wrapper-generation. Clearly this approach does not replace prior techniques of form-filling and wrapper generation,

particularly in the presence of domain knowledge or a corpus of examples. Preliminary results with the FormSolve prototype show that the high-level approach, along with our bytecode-based implementation, has promise in practice. A full-featured implementation of bytecode-based symbolic tracing in a state-of-the-art browser remains a major engineering challenge. In the future we will look at variants of our approach with symbolic tracing at the JavaScript level, working with the code on-the-fly. *Our benchmark generators and evaluation data are available on GitHub [16].*

References

1. Saswat Anand, Mayur Naik, Mary Jean Harrold, and Hongseok Yang. Automated concolic testing of smartphone apps. In *FSE*, 2012.
2. Artemis. <https://github.com/cs-au-dk/Artemis/tree/ff6ee7ee>, 2017.
3. Shay Artzi, Julian Dolby, Simon Holm Jensen, Anders Møller, and Frank Tip. A framework for automated testing of JavaScript Web applications. In *ICSE*, 2011.
4. Luciano Barbosa and Juliana Freire. An adaptive crawler for locating hidden-Web entry points. In *WWW*, 2007.
5. Luciano Barbosa and Juliana Freire. Siphoning hidden-Web data through keyword-based interfaces. *JIDM*, 2010.
6. Clark Barrett, Leonardo de Moura, Silvio Ranise, Aaron Stump, and Cesare Tinelli. The SMT-LIB initiative and the rise of SMT. In *HVC*, 2010.
7. Michael Benedikt, Tim Furche, Andreas Savvides, and Pierre Senellart. ProFoUnd: Program-analysis-based form understanding. In *WWW*, 2012. Demonstration.
8. Michael K. Bergman. The deep Web: Surfacing hidden value. *J. Electronic Publishing*, 7, 2001.
9. David R. Cok, David Déharbe, and Tjark Weber. The 2014 SMT competition. *Journal on Satisfiability, Boolean Modeling and Computation*, 9:207–242, 2014.
10. Crawljax, version 3.5. <http://crawljax.com/>, 2014.
11. Valter Crescenzi, Giansalvatore Mecca, and Paolo Merialdo. RoadRunner: Towards automatic data extraction from large Web sites. In *VLDB*, 2001.
12. Leonardo de Moura and Nikolaj Bjørner. Z3: an efficient SMT solver. In *TACAS*, 2008.
13. Leonardo de Moura and Nikolaj Bjørner. Satisfiability modulo theories: Introduction and applications. *Communications of the ACM*, 54(9):69–77, 2011.
14. Cristian Duda, Gianni Frey, Donald Kossmann, Reto Matter, and Chong Zhou. AJAX crawl: Making AJAX applications searchable. In *ICDE*, 2009.
15. Bruno Dutertre. Yices 2.2. In *CAV*, 2014.
16. FormSolve test data. <https://github.com/BenSpencer/FormSolve-Test-Data>.
17. Tim Furche, Georg Gottlob, Giovanni Grasso, Xiaonan Guo, Giorgio Orsi, and Christian Schallhart. OPAL: Automated form understanding for the deep Web. In *WWW*, 2012.
18. Tim Furche, Georg Gottlob, Giovanni Grasso, Xiaonan Guo, Giorgio Orsi, Christian Schallhart, and Cheng Wang. DIADEM: thousands of websites to a single database. *PVLDB*, 7(14):1845–1856, 2014.
19. Vijay Ganesh and David L. Dill. A decision procedure for bit-vectors and arrays. In *CAV*, 2007.
20. Patrice Godefroid, Nils Klarlund, and Koushik Sen. DART: directed automated random testing. In *PLDI*, 2005.
21. Bin He, Mitesh Patel, Zhen Zhang, and Kevin Chen-Chuan Chang. Accessing the deep Web. *Communications of the ACM*, 50(5):94–101, May 2007.
22. Gang Hu, Xinhao Yuan, Yang Tang, and Junfeng Yang. Efficiently, effectively detecting mobile app bugs with AppDoctor. In *EuroSys*, 2014.
23. Jalangi. <https://github.com/SRA-SiliconValley/jalangi/tree/22f130ec>, 2015.

24. Casper S. Jensen, Mukul R. Prasad, and Anders Møller. Automated testing with targeted event sequence generation. In *ISSTA*, July 2013.
25. Lu Jiang, Zhaohui Wu, Qian Feng, Jun Liu, and Qinghua Zheng. Efficient deep Web crawling using reinforcement learning. In *PAKDD*, 2010.
26. Gustavo Zanini Kantorski, Tiago Guimaraes Moraes, Viviane Pereira Moreira, and Carlos Alberto Heuser. Choosing values for text fields in Web forms. In *Advances in Databases and Information Systems*, 2013.
27. Daniel Kroening and Ofer Strichman. *Decision Procedures: An Algorithmic Point of View*. Springer, 2nd edition, 2016.
28. Juliano Palmieri Lage, Altigran S. da Silva, Paulo B. Golgher, and Alberto H. F. Laender. Automatic generation of agents for collecting hidden Web pages for data extraction. *Data Knowl. Eng.*, 49(2):177–196, May 2004.
29. Guodong Li, Esben Andreasen, and Indradeep Ghosh. SymJS: Automatic symbolic testing of JavaScript Web applications. In *FSE*, 2014.
30. Panagiotis Liakos, Alexandros Ntoulas, Alexandros Labrinidis, and Alex Delis. Focused crawling for the hidden Web. In *WWW*, 2016.
31. Tianyi Liang, Andrew Reynolds, Cesare Tinelli, Clark Barrett, and Morgan Deters. A DPLL(T) theory solver for a theory of strings and regular expressions. In *CAV*, 2014.
32. Jianguo Lu, Yan Wang, Jie Liang, Jessica Chen, and Jiming Liu. An approach to deep Web crawling by sampling. In *WI-IAT*, 2008.
33. Y. Lu, H. He, H. Zhao, W. Meng, and C. Yu. Annotating structured data of the deep Web. In *ICDE*, 2007.
34. Jayant Madhavan, David Ko, Lucja Kot, Vignesh Ganapathy, Alex Rasmussen, and Alon Halevy. Google’s deep-Web crawl. In *VLDB*, 2008.
35. Bruce McKenzie. The generation of strings from a CFG using a functional language, 1997. <https://ir.canterbury.ac.nz/handle/10092/11231>.
36. Ali Mesbah, Arie van Deursen, and Stefan Lenselink. Crawling Ajax-based Web applications through dynamic analysis of user interface state changes. *ACM Trans. Web*, 6(1):3:1–3:30, March 2012.
37. Hoa Nguyen, Thanh Nguyen, and Juliana Freire. Learning to extract form labels. *PVLDB*, 1(1):684–694, August 2008.
38. Aina Niemetz, Mathias Preiner, and Armin Biere. Boolector 2.0 system description. *Journal on Satisfiability, Boolean Modeling and Computation*, 9:53–58, 2015.
39. Alexandros Ntoulas, Petros Zerfos, and Junghoo Cho. Downloading textual hidden Web content through keyword queries. In *JCDL*, 2005.
40. Sriram Raghavan and Hector Garcia-Molina. Crawling the hidden Web. In *VLDB*, 2001.
41. Koushik Sen, Swaroop Kalasapur, Tasneem G. Brutch, and Simon Gibbs. Jalangi: a tool framework for concolic testing, selective record-replay, and dynamic analysis of JavaScript. In *FSE*, 2013.
42. Koushik Sen, Darko Marinov, and Gul Agha. CUTE: a concolic unit testing engine for C. In *FSE*, 2005.
43. Pierre Senellart, Avin Mittal, Daniel Muschick, Rémi Gilleron, and Marc Tommasi. Automatic wrapper induction from hidden-Web sources with domain knowledge. In *WIDM*, 2008.
44. Ben Spencer, Michael Benedikt, Anders Møller, and Franck van Breugel. ArtForm: A tool for exploring the codebase of form-based websites. In *ISSTA*, 2017.
45. Yan Wang, Jie Liang, and Jianguo Lu. Discover hidden Web properties by random walk on bipartite graph. *Information Retrieval*, 17(3), Jun 2014.
46. Wensheng Wu, AnHai Doan, Clement Yu, and Weiyi Meng. Modeling and extracting deep-Web query interfaces. In *Advances in Information and Intelligent Systems*. 2009.
47. Hongkun Zhao, Weiyi Meng, Zonghuan Wu, Vijay Raghavan, and Clement Yu. Fully automatic wrapper generation for search engines. In *WWW*, 2005.